

A taxonomy of sublinear multiple keyword pattern matching algorithms

B.W. Watson, G. Zwaan*

*Department of Mathematics and Computing Science, Eindhoven University of Technology,
P.O. Box 513, 5600 MB Eindhoven, Netherlands*

Received April 1995; revised February 1996

Communicated by M. Rem

Abstract

This article presents a taxonomy of sublinear keyword pattern matching algorithms related to the Boyer–Moore algorithm [3] and the Commentz-Walter algorithm [5, 6]. The taxonomy includes, amongst others, the multiple keyword generalization of the single keyword Boyer–Moore algorithm and an algorithm by Fan and Su [9, 10]. The corresponding precomputation algorithms are presented as well. The taxonomy is based on the idea of ordering algorithms according to their essential problem and algorithm details, and deriving all algorithms from a common starting point by successively adding these details in a correctness preserving way. This way of presentation not only provides a complete correctness argument of each algorithm, but also makes very clear what algorithms have in common (the details of their nearest common ancestor) and where they differ (the details added after their nearest common ancestor). Introduction of the notion of safe shift distances proves to be essential in the derivation and classification of the algorithms. Moreover, the article provides a common derivation for and a uniform presentation of the precomputation algorithms, not yet found in the literature.

1. Introduction

The keyword (or string) pattern matching problem can informally be described as the problem of finding all occurrences of keywords (strings) from a given set as substrings in a given (input) string. This problem is encountered in many areas and in several forms. In computing science, for instance, it plays a role in text search/analysis, lexical analysis, and data processing. In biology it is encountered in the analysis of, amongst others, DNA sequences. The problem can also be generalized to the matching of regular expressions, tree patterns, and graph patterns, none of which is treated in this article.

The keyword pattern matching problem has been extensively studied and a multitude of diverse solutions/algorithms exists. Single keyword algorithms are, for instance,

* Corresponding author. E-mail: wswinswan@win.tue.nl.

described by Knuth et al. [14] and Boyer and Moore [3]. Multiple keyword algorithms are described by Aho and Corasick [2], by Commentz-Walter [5, 6], and by Fan and Su [9, 10]. An overview of keyword pattern matching algorithms can be found in [1].

Due to the diversity of the algorithms and their descriptions – that tend to be rather involved and verbal – it is hard to get a good overview and to make a sound comparison between algorithms. In order to fulfill these needs a taxonomy of keyword pattern matching algorithms was presented by Watson and Zwaan [18–20]. Here, we focus our attention on a part of that taxonomy containing a family of multiple keyword pattern matching algorithms that have a matching time that may be sublinear in the length of the input string (taking the number of symbol comparisons as a measure of matching time). Amongst others, it comprises the multiple keyword generalization of the single keyword Boyer–Moore algorithm [3], the Commentz-Walter algorithm [5, 6], and the algorithm by Fan and Su [9, 10] (only after deriving this algorithm we found its description by Fan and Su). Both the Boyer–Moore and Commentz-Walter algorithms provided the inspiration for our derivations and classifying principle.

The main results of this article are comprised in the taxonomy graph shown in Fig. 1 (a more detailed description of the graph is found further on in this section). This taxonomy graph can be viewed as a table of contents to this article. It was obtained in order to meet the following goals:

- the systematic and formal derivation of the algorithms from a common starting point through a series of refinements to either algorithm or problem
- to factor out common portions of (the derivations of) well-known algorithms in order to facilitate the understanding of these algorithms and their comparison
- the presentation of the algorithms in a common framework to permit an easier comprehension of and a better comparison between the algorithms.

It is the first concise and systematic presentation of and comparison between the algorithms from the family considered here. Another taxonomy of (single keyword) pattern matching algorithms by Hume and Sunday [12] does not meet the goals set in this article since there any derivations and proofs of algorithms are missing.

Moreover, we show that all functions that need to be precomputed for the pattern matching algorithms of this family can in a simple way be expressed in a small set of base functions. The definitions of the base functions can all be written in forms satisfying one general pattern. From this general pattern a precomputation algorithm scheme has been derived that can be instantiated for each base function to yield a precomputation algorithm for that function (sometimes, the resulting algorithm can be simplified and/or some trivial post-processing has to be added to it). Hence, the precomputation algorithms can also be derived and presented in a uniform way. This includes formal derivations of the precomputation algorithms for the Boyer–Moore algorithm, the Commentz-Walter algorithm and the algorithm presented by Fan and Su not yet found in the literature.

The taxonomy given here differs from the corresponding part of the taxonomy in [19, 20] in that we here present the correct multiple keyword generalization of the Boyer–Moore algorithm and a common ancestor of this algorithm and the Commentz-

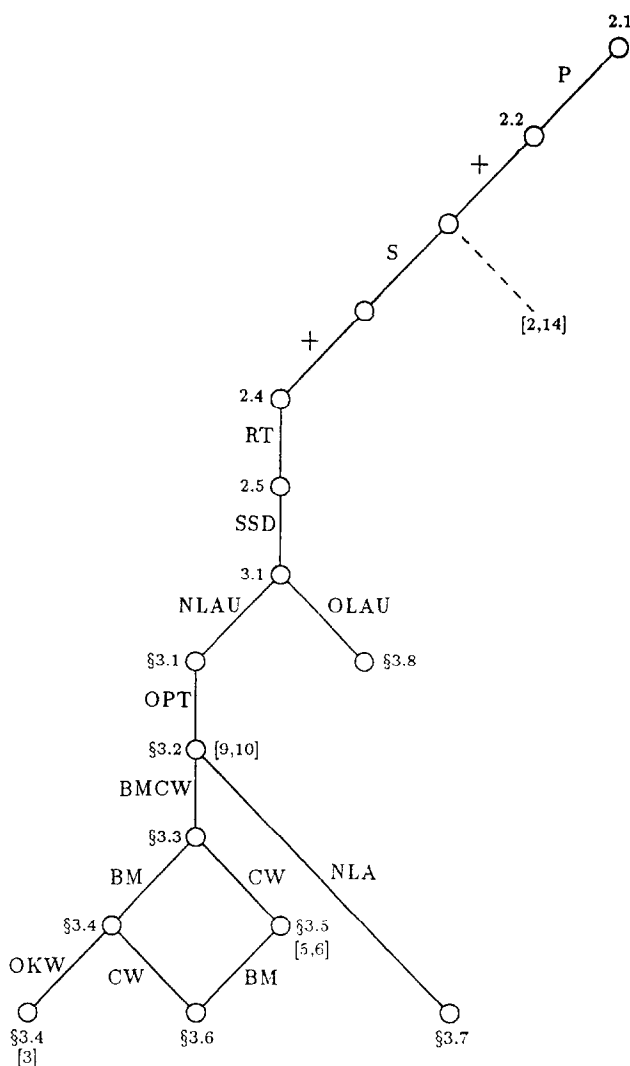


Fig. 1. Taxonomy graph of sublinear multiple keyword pattern matching algorithms. Each vertex corresponds to an algorithm. Vertices are labelled with either an algorithm number or a subsection number preceded by § referring to an algorithm or a subsection in this article. Some vertices are additionally labelled with literature references. Each edge corresponds to the addition of either a problem or algorithm detail (see Appendix C for a complete list of details and their descriptions). The sequence of edge labels that are encountered when going from the top vertex to another vertex forms a characterization of the algorithm corresponding to that vertex. For instance, the algorithm of the vertex labelled §3.5 (the Commentz-Walter algorithm) is characterized by (P+S+, RT, SSD, NLAU, OPT, BMCW, CW). The dashed edge leading to [2, 14] indicates the path to the rest of the taxonomy presented in [18–20] that contains the Knuth–Morris–Pratt and the Aho–Corasick algorithms. Since the algorithms in that part of the taxonomy are linear they fall outside the scope of this article and are omitted.

Walter algorithm. Because of this the derivations and the structure of the taxonomy have changed. Moreover, we present a completely new derivation of the precomputation algorithms.

An implementation as a C procedural library of almost all algorithms from [19] (called “Eindhoven Pattern Kit”) and an analysis of their performance is given in [17] (the implementation is in a somewhat rudimentary form being meant for the benchmarking only). Only the Commentz-Walter algorithm and a common descendant of both the Boyer–Moore and Commentz-Walter algorithm are implemented and discussed there. The performance results for these algorithms conform with the qualitative predictions made here. An implementation of almost all algorithms from this article and from [18, 19] as a C++ class library is called “SPARE Parts: A C++ toolkit for String PAttern REcognition” and is available at URL <ftp://ftp.win.tue.nl/pub/techreports/pi/pattm/spare/>. The implementation is entirely based on the abstract algorithms – in fact it is a systematic translation of them (this in contrast to, for instance, the implementations given in [12]). The SPARE Parts are fully described and documented in [18, Chapter 9].

1.1. Basic algorithm and derivation principles

The algorithms in this family traverse the input string in a direction that is opposite to the direction in which keyword symbols are matched to symbols in the input string. In this article we choose to inspect suffixes of prefixes of the input string both in order of increasing length. This algorithm will be the starting point of all further derivations. Choosing such a basic algorithm we have the possibility to attain matching times that are sublinear in the length of the input string (i.e. not all symbols of the input string are inspected). It is achieved by taking steps through the input string of which the length is determined by a shift function based on the information of the last matching attempt and possibly on additional information. The example in Fig. 2 illustrates this principle. Notice that not all symbols of the input string are scanned, although it is possible that some symbols of the input string are scanned more than once. The most simple shift function is the function that always yields 1. However, one can imagine larger shifts being possible. Ideally, such a shift would take us to the next occurrence of a match, but then calculating the value of the shift function is equivalent to the pattern matching problem itself. Therefore, we strive for shift functions that are easier to calculate and that do not exceed the ideal shift (called *safe shift functions*), i.e. we aim at approximations of the ideal shift from below. The ideal shift function is the minimum over a range characterized by some predicate. We derive various approximations from below by systematically weakening this predicate and derived predicates, by applying rules for minimum and maximum over disjunctive or conjunctive ranges, and by enlarging ranges. Considerations that play a role in these derivations are, for instance, whether or not to look ahead at symbols of the unscanned part of the input string, what information to use on the last scanned (non-matching) symbol, and the extent to which this information is coupled with the information on the

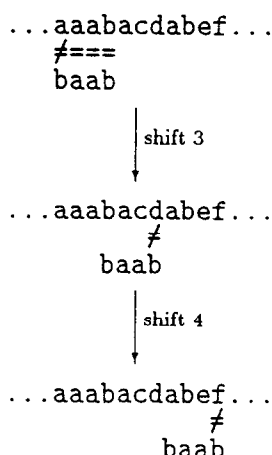


Fig. 2. Example of larger shift distances.

recognized suffix. Thus, we obtain several shift functions that meet the aforementioned requirements leading to an equal number of algorithms amongst which are the well-known Boyer–Moore and Commentz-Walter algorithms. The most simple weakening of the predicate is the weakening to the predicate true yielding the shift function that is always equal to 1. The reason to derive ever smaller shift functions is that a smaller shift function usually takes less precomputation time and less storage space (for instance, a one-dimensional table instead of a two-dimensional table).

The techniques we use to systematically derive shift functions from the ideal shift function enable us to clearly delineate the relations between the resulting shift functions and algorithms. Furthermore, they may be used in the derivation of yet other members of this family. The truly systematic approach, especially the predicate weakening technique, is not known from literature. Among other things we derive that the Commentz-Walter algorithm [5, 6] is not the multiple keyword generalization of the Boyer–Moore algorithm [3]. In fact, we show that the algorithms are incomparable (meaning that the shift distance in one algorithm is not always at least the shift distance in the other) and that they have a faster common ancestor that combines the properties of both. As it is, this common ancestor is derived first and subsequently the Boyer–Moore algorithm and the Commentz-Walter algorithm are derived from it. Both algorithms also have a common descendant (the algorithm that was incorrectly identified as the Boyer–Moore algorithm in [19, 20]). Furthermore, it is shown that the algorithm described by Fan and Su [9, 10] is an even faster ancestor of the common ancestor of the Boyer–Moore and Commentz-Walter algorithms.

1.2. The taxonomy

The algorithms are derived from a common starting point by successively adding either algorithm or problem details (see Appendix C for a complete list of details and

their descriptions). The only problem detail considered in this article is the restriction to the one keyword case. Among the algorithm details considered are restriction of nondeterminacy, introduction of the reverse trie, introduction of a shift function, and choice of a particular shift function. Each addition of an algorithm detail gives a new algorithm satisfying the same specification as the original algorithm; thus, the correctness of the algorithms is preserved. The sequence of details introduced in the derivation of an algorithm can be used to identify its similarities and its differences with other algorithms. These ordered detail sequences are used to identify the algorithms and to give a taxonomy of the algorithms in this family. The taxonomy is depicted as a graph in Fig. 1. Our method of developing a taxonomy was inspired by the method described by Jonkers [13]. There it is applied to develop a taxonomy of garbage collection algorithms. The method is also applied to attribute evaluation algorithms by Marcelis [15]. Other examples of algorithm taxonomies are found in [4, 7]. All algorithms are presented in a somewhat extended version of the guarded command language of Dijkstra [8] in order to avoid the peculiarities of any particular programming language. The algorithms use string type variables in order to abstract from any implementation detail like, for instance, indexing. Derivation of the algorithms is done a calculational way following Dijkstra's style of program derivation.

1.3. Overview

In Section 2 we give a formal definition of the pattern matching problem. From a trivial solution to this problem we derive, by addition of a number of algorithm details, an algorithm that is the starting point for the derivation of the algorithms in Section 3. In Section 3 we start by adding the algorithm detail that states that shifts larger than one may be possible. Addition of this program detail accounts for the possible sublinear matching time of all of the algorithms to be derived in this section. Subsequently, we derive by systematic approximation from below of the maximal safe shift distance the various algorithms of this family (in order of decreasing matching speed). The definitions of all functions introduced in Section 3 are rewritten in forms according to a general pattern in Section 4. From the general pattern a precomputation algorithm scheme is derived that can be instantiated for any of the functions. Section 5 contains the conclusions. Appendix A introduces a notation for quantifications and some of its properties. Appendix B contains definitions and properties used throughout this article. Appendix C contains a complete list of all algorithm and problem details and their descriptions.

2. The problem and some naive solutions

The keyword pattern matching problem is to find all occurrences of keywords from a set as substrings in an input string. Formally, given an alphabet V (a non-empty finite set of symbols), an input string $S \in V^*$, and a finite non-empty pattern set $P \subseteq V^*$,

establish (the notation used is described in Appendix A)¹

$$R: O = (\bigcup l, v, r : lvr = S : \{l\} \times (\{v\} \cap P) \times \{r\}).$$

A trivial (but unrealistic) solution to the problem is

Algorithm 2.1()

$$O := (\bigcup l, v, r : lvr = S : \{l\} \times (\{v\} \cap P) \times \{r\}) \\ \{R\}$$

The sequence of details describing this algorithm is the empty sequence (sequences of details are introduced in Section 1.2 and Fig. 1).

There are two basic directions in which to proceed while developing naive algorithms to solve this problem. Informally, a substring of S can be considered a “suffix of a prefix of S ” or a “prefix of a suffix of S ”. Only the first possibility is considered here, since the second possibility only leads to algorithms that are the mirror images of algorithms obtained by following the first possibility (basically, it amounts to reversing all strings in the problem). Moreover, this is the way that the Boyer–Moore, Commentz-Walter, and Fan and Su algorithms treat substrings of input string S .

Formally, we can consider “suffixes of prefixes of S ” as follows:

$$\begin{aligned} & (\bigcup l, v, r : lvr = S : \{l\} \times (\{v\} \cap P) \times \{r\}) \\ = & \{ \text{introduce } u : u = lv \} \\ & (\bigcup l, v, r, u : ur = S \wedge lv = u : \{l\} \times (\{v\} \cap P) \times \{r\}) \\ = & \{ l, v \text{ only occur in second range conjunct, so restrict their scope} \} \\ & (\bigcup u, r : ur = S : (\bigcup l, v : lv = u : \{l\} \times (\{v\} \cap P) \times \{r\})). \end{aligned}$$

A simple non-deterministic algorithm is obtained by applying “examine prefixes of a given string in any order” (algorithm detail (p)) to input string S . It results in²

Algorithm 2.2(p)

$$\begin{aligned} O &:= \emptyset; \\ \text{for } (u, r) : ur = S \rightarrow \\ & O := O \cup (\bigcup l, v : lv = u : \{l\} \times (\{v\} \cap P) \times \{r\}) \\ \text{rof } \{R\} \end{aligned}$$

The update of O (with another quantification) in the repetition of Algorithm 2.2 can be computed with another non-deterministic repetition. This inner repetition would

¹ Throughout this article we will adopt the convention that, unless stated otherwise, program variables and bound variables with names from the beginning of the Latin alphabet (i.e. a, b, c) will range over V , while variables with names from the end of the Latin alphabet (i.e. l, q, r, u, v, w) will range over V^* .

² In the following algorithms we use **for-rof** statements in order to express non-determinism. Statement **for** $x : P \rightarrow S$ **rof** amounts to executing statement list S once for each value of x that satisfies P initially, assuming only a finite number of such values exist. The order in which the values of x are chosen is arbitrary.

consider suffixes of u . Thus, by applying “examine suffixes of a given string in any order” (algorithm detail (s)) to string u we obtain algorithm

Algorithm 2.3(PS)

```

 $O := \emptyset;$ 
for  $(u, r) : ur = S \rightarrow$ 
  for  $(l, v) : lv = u \rightarrow$ 
     $O := O \cup \{l\} \times (\{v\} \cap P) \times \{r\}$ 
  rof
rof $\{R\}$ 

```

Algorithm 2.3 consists of two nested non-deterministic repetitions. In each case, the repetition can be made deterministic by considering prefixes (or suffixes as the case is) in increasing (called detail (+)) or decreasing (detail (−)) order of length. This gives two binary choices. Since the Boyer–Moore and Commentz-Walter algorithms examine string S from left to right and the patterns in P from right to left we focus our attention on (the operators $\uparrow, \downarrow, \vdash, \lceil$ are defined in Definition B.1; relation \leq_p is defined in Definition B.3):

Algorithm 2.4(P+S₊)

```

 $u, r := \varepsilon, S; \quad O := \{\varepsilon\} \times (\{\varepsilon\} \cap P) \times \{S\};$ 
 $\{\text{invariant: } O = (\bigcup x, y, z, : xyz = S \wedge xy \leq_p u : \{x\} \times (\{y\} \cap P) \times \{z\})\}$ 
do  $r \neq \varepsilon \rightarrow$ 
   $u, r := u(r \uparrow 1), r \downarrow 1;$ 
   $l, v := u, \varepsilon; \quad O := O \cup \{u\} \times (\{\varepsilon\} \cap P) \times \{r\};$ 
  do  $l \neq \varepsilon \rightarrow$ 
     $l, v := l \downarrow 1, (l \vdash 1)v;$ 
     $O := O \cup \{l\} \times (\{v\} \cap P) \times \{r\}$ 
  od
od $\{R\}$ 

```

Algorithm 2.4 has matching time $O(|S|^2)$, assuming that intersection with P is an $O(1)$ operation. We will now improve the matching time of this algorithm. Consider the set of suffixes of keywords $\text{su}ff(P)$ (the functions **su**ff and **pr**ef are defined in Definition B.2). A string w is an element of $\text{su}ff(P)$ if and only if it can be extended on the left to a pattern in P , i.e. $(\exists w' : w' \in V^* : w'w \in P)$. It follows that if $w \notin \text{su}ff(P)$ any extension of w on the left is not an element of $\text{su}ff(P)$ either. Consequently, the inner repetition in Algorithm 2.4 can terminate as soon as $(l \vdash 1)v \notin \text{su}ff(P)$ holds, since then all suffixes of u that are equal to or longer than $(l \vdash 1)v$ are not in $\text{su}ff(P)$ and hence not in P . The inner repetition guard is therefore strengthened to

$l \neq \varepsilon \text{ and } (l \vdash 1)v \in \text{su}ff(P).$

Observe that $v \in \text{su}ff(P)$ is now an invariant of the inner repetition. This invariant is initially established by the assignment $v := \varepsilon$ since $P \neq \emptyset$ and thus $\varepsilon \in \text{su}ff(P)$. Direct evaluation of $(l \upharpoonright 1)v \in \text{su}ff(P)$ is expensive. Therefore, it is done using the transition function τ_P of the *reverse trie* [11] corresponding to P where $\tau_P: \text{su}ff(P) \times V \rightarrow \text{su}ff(P) \cup \{\perp\}$ is defined for $w \in \text{su}ff(P)$ and $a \in V$ by

$$\tau_P(w, a) = \begin{cases} aw & \text{if } aw \in \text{su}ff(P), \\ \perp & \text{if } aw \notin \text{su}ff(P) \end{cases}$$

(algorithm detail (RT) (Reverse Trie)). Since we usually refer to the trie corresponding to P we will write τ instead of τ_P . Transition function τ can be computed beforehand. Its precomputation is discussed in Section 4.1. The guard becomes $l \neq \varepsilon$ **and** $\tau(v, l \upharpoonright 1) \neq \perp$ yielding algorithm:

Algorithm 2.5($P+S+RT$)

```

 $u, r := \varepsilon, S; O := \{\varepsilon\} \times (\{\varepsilon\} \cap P) \times \{S\};$ 
 $\{\text{invariant: } O = (\bigcup x, y, z : xyz = S \wedge xy \leq_p u : \{x\} \times (\{y\} \cap P) \times \{z\})\}$ 
do  $r \neq \varepsilon \rightarrow$ 
   $u, r := u(r \upharpoonright 1), r \downharpoonright 1;$ 
   $l, v := u, \varepsilon; O := O \cup \{u\} \times (\{\varepsilon\} \cap P) \times \{r\};$ 
  do  $l \neq \varepsilon$  and  $\tau(v, l \upharpoonright 1) \neq \perp \rightarrow$ 
     $l, v := l \upharpoonright 1, (l \upharpoonright 1)v;$ 
     $O := O \cup \{l\} \times (\{v\} \cap P) \times \{r\}$ 
  od
   $\{u = lv \wedge v \in \text{su}ff(P) \wedge (l = \varepsilon \text{ cor } (l \upharpoonright 1)v \notin \text{su}ff(P))\}$ 
od  $\{R\}$ 

```

This algorithm has $O(|S| \cdot (\text{MAX } p: p \in P : |p|))$ matching time. It will serve as a starting point for the derivation of all algorithms in the next section.

3. Sublinear pattern matching algorithms

In this section we derive sublinear pattern matching algorithms starting with Algorithm 2.5 by exploring the possibility of safely (without missing matches) making shifts of more than one symbol, i.e. replacing assignment $u, r := u(r \upharpoonright 1), r \downharpoonright 1$ by assignment $u, r := u(r \upharpoonright k), r \downharpoonright k$ for some k satisfying (the **MIN**-quantification is described in Appendix A)

$$1 \leq k \leq (\text{MIN } n: 1 \leq n \leq |r| \wedge \text{su}ff(u(r \upharpoonright n)) \cap P \neq \emptyset : n)$$

(algorithm detail (SSD) (safe shift distance)). The upper bound is the distance to the next match, the maximal safe shift distance. A number k satisfying this condition is called a *safe shift distance*. Since computing the upper bound on k is essentially the same as the problem we are trying to solve we aim at easier to compute approximations

Algorithm 3.1($P+S_+,RT,SSD$)

```

 $u, r := \varepsilon, S; \quad O := \{\varepsilon\} \times (\{\varepsilon\} \cap P) \times \{S\}; \quad l, v := \varepsilon, \varepsilon;$ 
 $\{\text{invariant: } O = (\bigcup x, y, z : xyz = S \wedge xy \leq_p u : \{x\} \times (\{y\} \cap P) \times \{z\})$ 
 $\quad \wedge u = lv \wedge v \in \text{succ}(P) \wedge (l = \varepsilon \text{ cor } (l \upharpoonright 1)v \notin \text{succ}(P))\}$ 
do  $r \neq \varepsilon \rightarrow$ 
 $u, r := u(r \upharpoonright k(l, v, r)), \quad r \downarrow k(l, v, r);$ 
 $l, v := u, \varepsilon; \quad O := O \cup \{u\} \times (\{\varepsilon\} \cap P) \times \{r\};$ 
do  $l \neq \varepsilon$  and  $\tau(v, l \upharpoonright 1) \neq \perp \rightarrow$ 
 $l, v := l \downarrow 1, (l \upharpoonright 1)v;$ 
 $O := O \cup \{l\} \times (\{v\} \cap P) \times \{r\}$ 
od
od  $\{R\}$ 

```

Fig. 3. Scheme for sublinear pattern matching algorithms.

from below of the upper bound. These are derived by systematically weakening the predicate $\text{succ}(u(r \upharpoonright n)) \cap P \neq \emptyset$ in the range of the upper bound, weakening resulting predicates, applying rules for minimum and maximum over a disjunctive or conjunctive range, and enlarging ranges. Considerations that play a role in these derivations are, for instance, whether or not to look ahead at symbols of the unscanned part of the input string, whether or not to use information on the last scanned symbol (usually a non-matching symbol), and the extent to which this information is coupled with the information on the recognized suffix. Thus, several algorithms are obtained amongst which the generalized version of the Boyer–Moore algorithm [3], the Commentz-Walter algorithm [5, 6], and the algorithm by Fan and Su [9, 10]. We derive ever smaller shift functions since the smaller the shift function the less precomputation time and storage space for the functions constituting the shift function is usually needed. For instance, the algorithm by Fan and Su [9, 10] is faster than the Commentz-Walter algorithm [5, 6], but it needs a two-dimensional table to store one of the functions constituting its shift function, whereas the Commentz-Walter algorithm only needs one-dimensional tables.

In the derivations we use part of the postcondition of the inner repetition in Algorithm 2.5 ($u = lv \wedge v \in \text{succ}(P)$). Adding $l, v := \varepsilon, \varepsilon$ to the initial assignments in Algorithm 2.5 turns $u = lv \wedge v \in \text{succ}(P)$ into an invariant of the outer repetition. Due to the dependence of the upper bound on l, v , and r we will aim at shift functions k that depend on l, v , and r and write $k(l, v, r)$. Hence, we arrive at the algorithm scheme in Fig. 3 for all algorithms to be derived in this section. Particular algorithms are obtained by substituting their shift functions for $k(l, v, r)$. Such a substitution may not always yield an algorithm that exactly corresponds with its original description in the literature; sometimes an additional transformation of the resulting algorithm is needed (for instance, a phase shift of the repetition; see [19] for a phase shifted version of the algorithm scheme). Conjunct $l = \varepsilon \text{ cor } (l \upharpoonright 1)v \notin \text{succ}(P)$ is added to the invariant in order to stress that provided l is non-empty symbol $l \upharpoonright 1$ is non-matching.

3.1. No lookahead at the unscanned part of the input string

In this section we derive an approximation from below of the upper bound on k that does not depend on r and that will be a starting point of most of our further derivations. In terms of algorithms this means that we refrain from looking ahead at the symbols of r , the yet unscanned part of the input string (algorithm detail (NLAU) (NO LOOKAHEAD at UNSCANNED part of the input string)). This is in accordance with most of the algorithms we are aiming at. One symbol lookahead at the unscanned part of the input string is discussed in Section 3.8. We derive

$$\begin{aligned}
 & (\text{MIN } n : 1 \leq n \leq |r| \wedge \text{suff}(u(r \upharpoonright n)) \cap P \neq \emptyset : n) \\
 \geq & \{n \leq |r|, r \upharpoonright n \in V^n, \text{ monotonicity of } \text{suff} \text{ and } \cap, \text{ enlarging range}\} \\
 & (\text{MIN } n : 1 \leq n \leq |r| \wedge \text{suff}(uV^n) \cap P \neq \emptyset : n) \\
 \geq & \{\text{enlarging range}\} \\
 & (\text{MIN } n : 1 \leq n \wedge \text{suff}(uV^n) \cap P \neq \emptyset : n).
 \end{aligned}$$

Since the last formula is to be the starting point of our further derivations we will from here on aim at shift functions k being dependent only on u , i.e. on l and v (remember $u = lv \wedge v \in \text{suff}(P)$). We will write $k(l, v)$ instead of $k(l, v, r)$.

3.2. Restriction to one symbol lookahead

In all derivations in this subsection and the following subsections we assume

$$u = lv \wedge v \in \text{suff}(P).$$

Restriction to one symbol lookahead ($l \upharpoonright 1$, the last symbol of u scanned in the inner loop) leads to the algorithm by Fan and Su [9, 10]. It is obtained by weakening the predicate in the approximation of the upperbound in subsection 3.1 in the following way:

$$\begin{aligned}
 & \text{suff}(uV^n) \cap P \neq \emptyset \\
 = & \{u = lv\} \\
 & \text{suff}(lvV^n) \cap P \neq \emptyset \\
 \Rightarrow & \{l = (l \upharpoonright 1)(l \upharpoonright 1), l \upharpoonright 1 \in V^*, \text{ monotonicity of } \text{suff} \text{ and } \cap\} \\
 & \text{suff}(V^*(l \upharpoonright 1)vV^n) \cap P \neq \emptyset \\
 = & \{\text{Property B.4}\} \\
 & V^*(l \upharpoonright 1)vV^n \cap V^*P \neq \emptyset \\
 = & \{l = \varepsilon : \text{Property B.5(i)}; l \neq \varepsilon : \text{Property B.5(ii)}\} \\
 & V^*(l \upharpoonright 1)vV^n \cap P \neq \emptyset \vee vV^n \cap V^*P \neq \emptyset.
 \end{aligned}$$

Notice that we have obtained a weaker predicate solely by discarding any information on $l \upharpoonright 1$. The only information on l that is still taken into account is $l \upharpoonright 1$, being either empty or consisting of one symbol. In the latter case we say to have one symbol lookahead. Observe that the symbol is the last symbol of u scanned in the inner loop

and that it is a non-matching symbol. After substituting the weaker predicate we obtain shift distance $k_{opt}(l, v)$ where $k_{opt} \in V^* \times \mathbf{su}\mathbf{ff}(P) \rightarrow \mathbb{N}$ is defined for $x \in V^*$ and $y \in \mathbf{su}\mathbf{ff}(P)$ by

$$k_{opt}(x, y) = (\mathbf{MIN} n : n \geq 1 \wedge (V^*(x \upharpoonright 1)yV^n \cap P \neq \emptyset \vee yV^n \cap V^*P \neq \emptyset) : n).$$

Having a disjunctive range (Property A.1) function k_{opt} can be expressed as follows:

$$k_{opt}(x, y) = \begin{cases} d_{opt}(x \upharpoonright 1, y) \mathbf{min} d_{sp}(y), & x \neq \varepsilon, \\ d_i(y) \mathbf{min} d_{sp}(y) & x = \varepsilon \end{cases}$$

where $d_{opt} \in V \times \mathbf{su}\mathbf{ff}(P) \rightarrow \mathbb{N}$ is defined for $a \in V$ and $y \in \mathbf{su}\mathbf{ff}(P)$ by

$$d_{opt}(a, y) = (\mathbf{MIN} n : n \geq 1 \wedge V^*ayV^n \cap P \neq \emptyset : n),$$

$d_{sp} \in \mathbf{su}\mathbf{ff}(P) \rightarrow \mathbb{N}$ is defined for $y \in \mathbf{su}\mathbf{ff}(P)$ by

$$d_{sp}(y) = (\mathbf{MIN} n : n \geq 1 \wedge yV^n \cap V^*P \neq \emptyset : n)$$

(cf. [5, 6, function d_2]), and $d_i \in \mathbf{su}\mathbf{ff}(P) \rightarrow \mathbb{N}$ is defined for $y \in \mathbf{su}\mathbf{ff}(P)$ by

$$d_i(y) = (\mathbf{MIN} n : n \geq 1 \wedge V^*yV^n \cap P \neq \emptyset : n)$$

(cf. [5, 6, function d_1]). Functions d_{opt} and d_i account for occurrences of ay and y , respectively, within some keyword (i.e. as infix of some keyword), whereas function d_{sp} accounts for occurrences of suffixes of y as proper prefixes of some keyword. Precomputation of d_{opt} , d_{sp} , and d_i is discussed in Section 4.2. Calculating the shift distance in this way is referred to as algorithm detail (OPT) (OPTimal for one symbol lookahead) and results in algorithm $(P_+S_+,RT,SSD,NLAU,OPT)$. We arrived at this algorithm not knowing it had already been described by Fan and Su in [9, 10]. From their informal description it undoubtedly follows that they describe the same algorithm though their formal treatment of the algorithm and, especially, the precomputation is rather involved. Finally, notice that to store function d_{opt} one needs a two-dimensional table, whereas functions d_i and d_{sp} only need one-dimensional tables. In the following sections we derive shift functions smaller than k_{opt} that are expressed solely in functions needing one-dimensional tables for storage.

3.3. Lookahead symbol is mismatching

We derive an approximation from below of d_{opt} that yields an algorithm that is the common ancestor of the multiple keyword generalization of the Boyer–Moore algorithm [3] and the Commentz-Walter algorithm [5, 6]. Essentially, the resulting shift function is not based on the identity of the lookahead symbol $l \upharpoonright 1$ but only uses the fact that the lookahead symbol is mismatching, as is done in the Boyer–Moore shift function. In this way one might say that the recognized suffix and the (mismatching) lookahead symbol have to some extent been decoupled.

We start by weakening the range predicate from d_{opt} . Let $a \in V$ and $y \in \mathbf{suff}(P)$ such that $ay \notin \mathbf{suff}(P)$. We derive

$$\begin{aligned}
 & V^*ayV^n \cap P \neq \emptyset \\
 = & \{y \in V^{|y|}, \text{ monotonicity of } \cap\} \\
 & V^*aV^{|y|+n} \cap P \neq \emptyset \wedge V^*ayV^n \cap P \neq \emptyset \\
 \Rightarrow & \left\{ \begin{array}{l} ay \notin \mathbf{suff}(P), \text{ so } a \in \{b \mid b \in V \wedge by \notin \mathbf{suff}(P)\}, \\ \text{definition } MS \text{ (after derivation)} \end{array} \right\} \\
 & V^*aV^{|y|+n} \cap P \neq \emptyset \wedge V^*(V \setminus MS(y))yV^n \cap P \neq \emptyset
 \end{aligned}$$

where $MS \in \mathbf{suff}(P) \rightarrow \mathcal{P}(V)$ is defined for $z \in \mathbf{suff}(P)$ by

$$MS(z) = \{b \mid b \in V \wedge bz \in \mathbf{suff}(P)\}.$$

The first conjunct will lead to a shift component based on the identity of the lookahead symbol that is identical to a component of the Commentz-Walter shift function. The second conjunct will lead to a shift component – based on the recognized suffix and the fact that the lookahead symbol is mismatching – that is identical to a component of the Boyer–Moore shift function. Assuming $l \neq \varepsilon$ and $(l \uparrow 1)v \notin \mathbf{suff}(P)$ we derive

$$\begin{aligned}
 & d_{opt}(l \uparrow 1, v) \\
 \geq & \{\text{definition } d_{opt}, \text{ preceding derivation, enlarging range}\} \\
 & (\mathbf{MIN} \, n : n \geq 1 \wedge V^*(l \uparrow 1)V^{|v|+n} \cap P \neq \emptyset \\
 & \wedge V^*(V \setminus MS(v))vV^n \cap P \neq \emptyset : n) \\
 \geq & \{\mathbf{MIN} \text{ with conjunctive range}\} \\
 & (\mathbf{MIN} \, n : n \geq 1 \wedge V^*(l \uparrow 1)V^{|v|+n} \cap P \neq \emptyset : n) \\
 & \mathbf{max} (\mathbf{MIN} \, n : n \geq 1 \wedge V^*(V \setminus MS(v))vV^n \cap P \neq \emptyset : n) \\
 = & \left\{ \begin{array}{l} \text{change of bound variable: } m = |v| + n; \\ \text{definition } d_{vi} \text{ (after derivation)} \end{array} \right\} \\
 & (\mathbf{MIN} \, m : m \geq |v| + 1 \wedge V^*(l \uparrow 1)V^m \cap P \neq \emptyset : m - |v|) \mathbf{max} \, d_{vi}(v) \\
 \geq & \{\text{enlarging range}\} \\
 & (\mathbf{MIN} \, m : m \geq 1 \wedge V^*(l \uparrow 1)V^m \cap P \neq \emptyset : m - |v|) \mathbf{max} \, d_{vi}(v) \\
 = & \{l \neq \varepsilon, \text{ definition of } char_{cw} \text{ (after derivation)}\} \\
 & char_{cw}(l \uparrow 1, |v|) \mathbf{max} \, d_{vi}(v)
 \end{aligned}$$

where $d_{vi} \in \mathbf{suff}(P) \rightarrow \mathbb{N}$ is defined for $y \in \mathbf{suff}(P)$ by

$$d_{vi}(y) = (\mathbf{MIN} \, n : n \geq 1 \wedge V^*(V \setminus MS(y))yV^n \cap P \neq \emptyset : n)$$

and $char_{cw} \in V \times \mathbb{N} \rightarrow \mathbb{N}$ is defined for $a \in V$ and $z \in \mathbb{N}$ by

$$char_{cw}(a, z) = (\mathbf{MIN} \, n : n \geq 1 \wedge V^*aV^n \cap P \neq \emptyset : n - z).$$

This results in shift distance $k_{bmcw}(l, v)$ where $k_{bmcw} \in V^* \times \mathbf{su}\mathbf{ff}(P) \rightarrow \mathbb{N}$ is defined for $x \in V^*$ and $y \in \mathbf{su}\mathbf{ff}(P)$ by

$$k_{bmcw}(x, y) = \begin{cases} (\mathit{char}_{cw}(x \upharpoonright 1, |y|) \mathbf{max} d_{vi}(y)) \mathbf{min} d_{sp}(y) & \text{if } x \in V^+, \\ d_i(y) \mathbf{min} d_{sp}(y) & \text{if } x = \varepsilon. \end{cases}$$

Function char_{cw} can be expressed for $a \in V$ and $z \in \mathbb{N}$ as

$$\mathit{char}_{cw}(a, z) = \begin{cases} +\infty & \text{if } \overline{\mathit{char}}_{cw}(a) = +\infty, \\ \overline{\mathit{char}}_{cw}(a) - z & \text{if } \overline{\mathit{char}}_{cw}(a) \neq +\infty, \end{cases}$$

where function $\overline{\mathit{char}}_{cw} \in V \rightarrow \mathbb{N}$ is defined for $a \in V$ by

$$\overline{\mathit{char}}_{cw}(a) = (\mathbf{MIN} n : n \geq 1 \wedge V^* a V^n \cap P \neq \emptyset : n).$$

From the definition $\overline{\mathit{char}}_{cw}$ it immediately follows that its computation can be interwoven with the precomputation of the reverse trie τ that is discussed in Section 4.1. Precomputation of d_{vi} , d_{sp} , and d_i is discussed in Section 4.2.

Approximation from below of k_{opt} by k_{bmcw} is referred to as algorithm detail (BMCW). We chose this name to reflect that essential ideas from both the Boyer–Moore and Commentz-Walter algorithms are introduced. In the next two subsections these algorithms are derived from the algorithm presented in this subsection and characterized by detail sequence (P+S+,RT,SSD,NLAU,OPT,BMCW). Notice that for $x \in V^*$ and $y \in \mathbf{su}\mathbf{ff}(P)$ we have

$$k_{opt}(x, y) \geq k_{bmcw}(x, y).$$

3.4. The multiple keyword Boyer–Moore algorithm

We proceed by deriving the multiple keyword generalization of the Boyer–Moore algorithm [3] from the algorithm in Section 3.3. It only differs from the algorithm there in the way the lookahead symbol is taken into account. Assuming $l \neq \varepsilon$ we derive

$$\begin{aligned} & \mathit{char}_{cw}(l \upharpoonright 1, |v|) \\ = & \{ \text{definition } \mathit{char}_{cw} \} \\ & (\mathbf{MIN} n : n \geq 1 \wedge V^*(l \upharpoonright 1) V^n \cap P \neq \emptyset : n - |v|) \\ \geq & \{ V^*(l \upharpoonright 1) V^n \cap P \neq \emptyset \Rightarrow V^*(l \upharpoonright 1) V^n \cap V^* P \neq \emptyset, \text{ enlarging range} \} \\ & (\mathbf{MIN} n : n \geq 1 \wedge V^*(l \upharpoonright 1) V^n \cap V^* P \neq \emptyset : n - |v|) \\ = & \left\{ \begin{array}{l} P \neq \emptyset, \text{ take } p \in P \text{ then } V^*(l \upharpoonright 1) V^{|p|+1} \cap V^* P \neq \emptyset, \\ |p| + 1 \geq 1, \text{ non-empty range} \end{array} \right\} \\ & (\mathbf{MIN} n : n \geq 1 \wedge V^*(l \upharpoonright 1) V^n \cap V^* P \neq \emptyset : n) - |v| \\ = & \{ l \neq \varepsilon, \text{ definition of } \mathit{char}_{bm} \text{ (after derivation)} \} \\ & \mathit{char}_{bm}(l \upharpoonright 1) - |v| \end{aligned}$$

where $\text{char}_{bm} \in V \rightarrow \mathbb{N}$ is defined for $a \in V$ by

$$\text{char}_{bm}(a) = (\text{MIN } n : n \geq 1 \wedge V^* a V^n \cap V^* P \neq \emptyset : n).$$

It results in shift distance $k_{bm}(l, v)$ where $k_{bm} \in V^* \times \text{suff}(P) \rightarrow \mathbb{N}$ is defined for $x \in V^*$ and $y \in \text{suff}(P)$ by

$$k_{bm}(x, y) = \begin{cases} ((\text{char}_{bm}(x \upharpoonright 1) - |y|) \max d_{vi}(y)) \min d_{sp}(y) & \text{if } x \in V^+, \\ d_i(y) \min d_{sp}(y) & \text{if } x = \varepsilon. \end{cases}$$

Precomputation of d_{vi} , d_{sp} , and d_i is discussed in Section 4.2. Function char_{bm} can be expressed in terms of function $\overline{\text{char}}_{cw}$ (introduced in Section 3.3) as the following derivation shows. For $a \in V$ we derive

$$\begin{aligned} & \text{char}_{bm}(a) \\ = & \{ \text{definition } \text{char}_{bm} \} \\ & (\text{MIN } n : n \geq 1 \wedge V^* a V^n \cap V^* P \neq \emptyset : n) \\ = & \{ \text{Property B.5(iii)} \} \\ & (\text{MIN } n : n \geq 1 \wedge (V^* a V^n \cap P \neq \emptyset \vee a V^n \cap V^+ P \neq \emptyset) : n) \\ = & \{ \text{disjunctive range, definition } \overline{\text{char}}_{cw} \} \\ & \overline{\text{char}}_{cw}(a) \min (\text{MIN } n : n \geq 1 \wedge a V^n \cap V^+ P \neq \emptyset : n) \\ = & \{ a V^n \cap V^+ P \neq \emptyset \equiv V^n \cap V^* P \neq \emptyset \} \\ & \overline{\text{char}}_{cw}(a) \min (\text{MIN } n : n \geq 1 \wedge V^n \cap V^* P \neq \emptyset : n) \\ = & \{ \text{definition } m_P \text{ (after derivation)} \} \\ & \overline{\text{char}}_{cw}(a) \min m_P \end{aligned}$$

where

$$m_P = \begin{cases} 1 & \text{if } \varepsilon \in P, \\ (\text{MIN } p : p \in P : |p|) & \text{if } \varepsilon \notin P. \end{cases}$$

Approximating $k_{bm_{cw}}$ from below by k_{bm} is referred to as algorithm detail (BM). It results in the multiple keyword generalization of the regular Boyer–Moore algorithm [3]. The algorithm is characterized by detail sequence $(P_+, S_+, RT, SSD, NLAU, OPT, BMCW, BM)$. The regular Boyer–Moore algorithm can be obtained by restricting P to one keyword (problem detail (OKW) (one keyword)). Notice that for $x \in V^*$ and $y \in \text{suff}(P)$ we have

$$k_{bm_{cw}}(x, y) \geq k_{bm}(x, y).$$

Inequality can only occur if the lookahead symbol does not occur in any keyword except as the rightmost symbol, since in that case char_{cw} yields $+\infty$ (being the minimum over an empty range (see Appendix A)), whereas char_{bm} yields the finite value m_P .

The formula for the Boyer–Moore shift function given here differs from the ones given in [1, 3]. We will show that all formulas are equivalent. We have for all $a \in V$

$$\text{char}_{bm}(a) \leq m_P.$$

Since for all $y \in \text{suff}(P)$ we have $(\forall n : 1 \leq n < m_P - |y| : yV^n \cap V^*P = \emptyset)$, it follows that for all $y \in \text{suff}(P)$

$$m_P - |y| \leq d_{sp}(y).$$

Finally, we derive for $x \in V^+$ and $y \in \text{suff}(P)$

$$\begin{aligned} & k_{bm}(x, y) \\ = & \{ \text{definition } k_{bm} \} \\ & ((\text{char}_{bm}(x \upharpoonright 1) - |y|) \max d_{vi}(y)) \min d_{sp}(y) \\ = & \{ d_{sp}(y) \geq m_P - |y| \geq \text{char}_{bm}(x \upharpoonright 1) - |y| \} \\ & ((\text{char}_{bm}(x \upharpoonright 1) - |y|) \max d_{vi}(y)) \min ((\text{char}_{bm}(x \upharpoonright 1) - |y|) \max d_{sp}(y)) \\ = & \{ \text{distributivity} \} \\ & (\text{char}_{bm}(x \upharpoonright 1) - |y|) \max (d_{vi}(y) \min d_{sp}(y)). \end{aligned}$$

The last formula in the preceding derivation coincides with the ones in [1, 3].

3.5. The Commentz-Walter algorithm

Instead of approximating char_{cw} in k_{bmcw} from below by char_{bm} we now approximate d_{vi} in k_{bmcw} from below by d_i . This results in the Commentz-Walter algorithm [5, 6]. We derive

$$\begin{aligned} & d_{vi}(v) \\ = & \{ \text{definition } d_{vi} \} \\ & (\text{MIN } n : n \geq 1 \wedge V^*(V \setminus MS(v))vV^n \cap P \neq \emptyset : n) \\ \geq & \left\{ \begin{array}{l} V^*(V \setminus MS(v))vV^n \cap P \neq \emptyset \Rightarrow V^*vV^n \cap P \neq \emptyset, \\ \text{enlarging range} \end{array} \right\} \\ & (\text{MIN } n : n \geq 1 \wedge V^*vV^n \cap P \neq \emptyset : n) \\ = & \{ \text{definition } d_i \} \\ & d_i(v). \end{aligned}$$

This results in shift distance $k_{cw}(l, v)$ where $k_{cw} \in V^* \times \text{suff}(P) \rightarrow \mathbb{N}$ is defined for $x \in V^*$ and $y \in \text{suff}(P)$ by

$$k_{cw}(x, y) = \begin{cases} (\text{char}_{cw}(x \upharpoonright 1, |y|) \max d_i(y)) \min d_{sp}(y) & \text{if } x \in V^+, \\ d_i(y) \min d_{sp}(y) & \text{if } x = \varepsilon. \end{cases}$$

Function $char_{cw}$ has been introduced in Section 3.3. Precomputation of d_i and d_{sp} is discussed in Section 4.2.

Approximating k_{bmcw} from below by k_{cw} is referred to as algorithm detail (cw). It results in the Commentz-Walter algorithm [5, 6] that is characterized by detail sequence $(P_+, S_+, RT, SSD, NLAU, OPT, BMCW, CW)$. Notice that for $x \in V^*$ and $y \in \text{succ}(P)$ we have

$$k_{bmcw}(x, y) \geq k_{cw}(x, y).$$

Such a comparison cannot be made between k_{hm} and k_{cw} as the following example shows.

Example 3.1. Let $V = \{a, b, c, d\}$, $P = \{cababa\}$, and $x \in V^*$. Shift functions k_{hm} and k_{cw} are incomparable since

$$\begin{aligned} k_{opt}(xd, a) &= d_{opt}(d, a) \min d_{sp}(a) \\ &= +\infty \min 6 = 6 \\ k_{bmcw}(xd, a) &= (char_{cw}(d, 1) \max d_{vi}(a)) \min d_{sp}(a) \\ &= (+\infty \max 4) \min 6 = 6 \\ k_{hm}(xd, a) &= ((char_{hm}(d) - 1) \max d_{vi}(a)) \min d_{sp}(a) \\ &= ((6 - 1) \max 4) \min 6 = 5 \\ k_{cw}(xd, a) &= (char_{cw}(d, 1) \max d_i(a)) \min d_{sp}(a) \\ &= (+\infty \max 2) \min 6 = 6 \end{aligned}$$

and

$$\begin{aligned} k_{opt}(xa, a) &= d_{opt}(a, a) \min d_{sp}(a) \\ &= +\infty \min 6 = 6 \\ k_{bmcw}(xa, a) &= (char_{cw}(a, 1) \max d_{vi}(a)) \min d_{sp}(a) \\ &= ((2 - 1) \max 4) \min 6 = 4 \\ k_{hm}(xa, a) &= ((char_{hm}(a) - 1) \max d_{vi}(a)) \min d_{sp}(a) \\ &= ((2 - 1) \max 4) \min 6 = 4 \\ k_{cw}(xa, a) &= (char_{cw}(a, 1) \max d_i(a)) \min d_{sp}(a) \\ &= ((2 - 1) \max 2) \min 6 = 2. \end{aligned}$$

It also follows that in some cases k_{bmcw} is smaller than k_{opt} and that in some cases k_{hm} or k_{cw} is smaller than k_{bmcw} .

It is not possible that both $k_{bmcw}(x, y) > k_{hm}(x, y)$ and $k_{bmcw}(x, y) > k_{cw}(x, y)$ hold for some $x \in V^+$ and $y \in \text{succ}(P)$ since the first inequality implies $char_{cw}(x \uparrow 1, |y|) = +\infty$ and this in its turn implies $k_{bmcw}(x, y) = d_{sp}(y) = k_{cw}(x, y)$.

3.6. Complete decoupling of recognized suffix and lookahead symbol

The derivations in the previous subsections effect an ever stronger decoupling of the recognized suffix v and the lookahead symbol $l \uparrow 1$ in the subsequent shift functions. By approximating d_{vi} in k_{bm} from below by d_i or $char_{cw}$ in k_{cw} by $char_{bm}$ (or both in k_{bmcw}) we obtain a complete decoupling. It results in shift distance $k_{dsl}(l, v)$ where $k_{dsl} \in V^* \times \mathbf{suff}(P) \rightarrow \mathbb{N}$ is defined for $x \in V^*$ and $y \in \mathbf{suff}(P)$ by

$$k_{dsl}(x, y) = \begin{cases} ((char_{bm}(x \uparrow 1) - |y|) \max d_i(y)) \min d_{sp}(y) & \text{if } x \in V^+, \\ d_i(y) \min d_{sp}(y) & \text{if } x = \varepsilon. \end{cases}$$

Function $char_{bm}$ has been introduced in Section 3.4. Precomputation of d_i and d_{sp} is discussed in Section 4.2. The algorithm can be characterized both by detail sequence $(P_+S_+, RT, SSD, NLAU, OPT, BMCW, BM, CW)$ and detail sequence $(P_+S_+, RT, SSD, NLAU, OPT, BMCW, CW, BM)$.

3.7. Discarding the lookahead symbol

We weaken the predicate in the range of k_{opt} by weakening its first disjunct to $V^*yV^n \cap P \neq \emptyset$ due to $V^*(x \uparrow 1) \subseteq V^*$ and the monotonicity of \cap . This weakening step is referred to as discarding the lookahead symbol $l \uparrow 1$. The shift distance corresponding to this weakening is $k_{nla}(v)$ where $k_{nla} \in \mathbf{suff}(P) \rightarrow \mathbb{N}$ is defined for $y \in \mathbf{suff}(P)$ by

$$k_{nla}(y) = d_i(y) \min d_{sp}(y).$$

Notice that this shift function can also be viewed as an approximation from below of k_{dsl} . Precomputation of d_i and d_{sp} is discussed in Section 4.2. Approximating k_{opt} from below by k_{nla} is referred to as algorithm detail (NLA) (NO Lookahead at mismatching symbol) and results in the algorithm characterized by detail sequence $(P_+S_+, RT, SSD, NLAU, OPT, NLA)$.

3.8. One symbol lookahead at the unscanned part of the input string

In this section we consider looking ahead at the first symbol of the unscanned part r of the input string. The first symbol of r will be taken into account independently of the other available information. In this way we obtain stronger variants of all of the shift functions derived thus far. Assuming $r \neq \varepsilon$ we derive

$$\begin{aligned} & (\mathbf{MIN} n : 1 \leq n \leq |r| \wedge \mathbf{suff}(u(r \uparrow n)) \cap P \neq \emptyset : n) \\ &= \{1 \leq n \leq |r|, r \uparrow n = (r \uparrow 1)((r \downarrow 1) \uparrow (n-1))\} \\ & (\mathbf{MIN} n : 1 \leq n \leq |r| \wedge \mathbf{suff}(u(r \uparrow 1)((r \downarrow 1) \uparrow (n-1))) \cap P \neq \emptyset : n) \\ &\geq \left\{ \begin{array}{l} 1 \leq n \leq |r|, (r \downarrow 1) \uparrow (n-1) \in V^{n-1}, \\ \text{monotonicity of } \mathbf{suff} \text{ and } \cap, \text{ enlarging range} \end{array} \right\} \end{aligned}$$

$$\begin{aligned}
& (\text{MIN } n : 1 \leq n \leq |r| \wedge \text{succ}(u(r \uparrow 1)V^{n-1}) \cap P \neq \emptyset : n) \\
& \geq \{u \in V^*, \text{ monotonicity of succ and } \cap, \text{ enlarging range}\} \\
& (\text{MIN } n : 1 \leq n \leq |r| \wedge \text{succ}(V^*(r \uparrow 1)V^{n-1}) \cap P \neq \emptyset : n) \\
& \geq \{\text{enlarging range, changing bound variable: } m = n - 1\} \\
& (\text{MIN } m : 0 \leq m \wedge \text{succ}(V^*(r \uparrow 1)V^m) \cap P \neq \emptyset : m + 1) \\
& = \left\{ \begin{array}{l} \text{Property B.4, } P \neq \emptyset, \text{ take } p \in P \text{ then} \\ V^*(r \uparrow 1)V^{|p|} \cap V^*p \neq \emptyset \text{ and } |p| \geq 0, \text{ non-empty range} \end{array} \right\} \\
& (\text{MIN } m : 0 \leq m \wedge V^*(r \uparrow 1)V^m \cap V^*P \neq \emptyset : m) + 1 \\
& = \{r \neq \varepsilon, \text{ definition } \text{char}_{la}(\text{after derivation})\} \\
& \text{char}_{la}(r \uparrow 1) + 1
\end{aligned}$$

where $\text{char}_{la} \in V \rightarrow \mathbb{N}$ is defined for $a \in V$ by

$$\text{char}_{la}(a) = (\text{MIN } n : 0 \leq n \wedge V^*aV^n \cap V^*P \neq \emptyset : n).$$

Function char_{la} can be expressed in terms of function char_{bm} (and hence in terms of $\overline{\text{char}}_{cw}$) as the following derivation shows. For $a \in V$ we derive

$$\begin{aligned}
& \text{char}_{la}(a) \\
& = \{\text{definition } \text{char}_{la}\} \\
& (\text{MIN } n : 0 \leq n \wedge V^*aV^n \cap V^*P \neq \emptyset : n) \\
& = \{\text{range split: } n \geq 1 \vee n = 0, \text{ definition } \text{char}_{bm}\} \\
& \text{char}_{bm}(a) \text{min} (\text{MIN } n : n = 0 \wedge V^*a \cap V^*P \neq \emptyset : 0) \\
& = \{\text{Property B.5(ii)}\} \\
& \text{char}_{bm}(a) \text{min} (\text{MIN } n : n = 0 \wedge (a \in \text{succ}(P) \vee \varepsilon \in P) : 0).
\end{aligned}$$

Let $M(u, r)$ denote the first expression in the first derivation of this subsection as well as the first expression in the derivation in Section 3.1, and let $N(u)$ denote the last expression in the derivation in Section 3.1. We have

$$\begin{aligned}
& M(u, r) \\
& = \{\text{property max}\} \\
& M(u, r) \text{max } M(u, r) \\
& \geq \{\text{derivation in Section 3.1, first derivation in this section}\} \\
& N(u) \text{max } (\text{char}_{la}(r \uparrow 1) + 1).
\end{aligned}$$

Since all shift functions derived in the previous subsections are approximations from below of $N(u)$ the preceding derivations shows that they all may be extended with $\text{max}(\text{char}_{la}(r \uparrow 1) + 1)$ to form a class of stronger shift functions of signature $k(l, v, r)$

(algorithm detail (OLAU) (one symbol LookAhead at unscanned part of the input string)). The first derivation in this subsection shows that it is also possible to couple the information on $r \uparrow 1$ with the information on l and v ($u = lv$). We will not pursue that direction any further in this article.

4. Precomputation

In this section we derive algorithms for the precomputation of the functions used in the pattern matching algorithms in Sections 2 and 3. The algorithms are correct due to their formal derivation. This cannot always be said about the algorithms found in the literature, mostly due to the absence of any formal derivation (see, for instance, the single keyword Boyer–Moore precomputation algorithms given in [3, 14, 16], where each article shows the preceding article to give an incorrect precomputation algorithm). Moreover, we give the first formal derivation of the precomputation algorithms for the family of sublinear pattern matching algorithms. They can, amongst others, be specialized to a correct precomputation algorithm for the single keyword Boyer–Moore algorithm. In fact, we show that the definition of all d -functions introduced in Section 3 can be rewritten into a form in accordance with one general pattern. Subsequently, a general precomputation algorithm scheme for this general pattern is derived that can be instantiated for every d -function.

4.1. Precomputation of τ_P

The transition function $\tau_P \in \text{succ}(P) \times V \rightarrow (\text{succ}(P) \cup \{\perp\})$ of the reverse trie corresponding to P is defined for $u \in \text{succ}(P)$ and $a \in V$ by

$$\tau_P(u, a) = \begin{cases} au & \text{if } au \in \text{succ}(P), \\ \perp & \text{if } au \notin \text{succ}(P). \end{cases}$$

Since **succ** is idempotent and the definition of τ_P only depends on **succ**(P), we have $\tau_P = \tau_{\text{succ}(P)}$. Set P being non-empty we have $\text{succ}(P) = \{\varepsilon\} \cup \text{succ}(P)$ and $\tau_{\text{succ}(P)} = \tau_{\{\varepsilon\} \cup \text{succ}(P)}$. These observations lead to the following algorithm (cf. [2, Section 3, Algorithm 2]) to compute τ_P in which variable τ is used to calculate and store τ_P thereby viewing τ as a set of ordered pairs (the usual notion of a function) and abbreviating statements like $\tau := \tau + \{(x, a), y\}$ to $\tau(x, a) := y$:

```

 $\tau := \emptyset; \{ \tau = \tau_\emptyset \}$ 
for  $a : a \in V \rightarrow \tau(\varepsilon, a) := \perp$  rof;  $\{ \tau = \tau_{\{\varepsilon\}} \}$ 
 $P_d, P_r := \emptyset, P;$ 
 $\{ \text{invariant: } P_d \cup P_r = P \wedge P_d \cap P_r = \emptyset \wedge \tau = \tau_{\{\varepsilon\} \cup \text{succ}(P_d)} \}$ 
do  $P_r \neq \emptyset \rightarrow$ 
   $p : p \in P_r; u, v := p, \varepsilon;$ 
   $\{ \text{invariant: } uv = p \wedge \tau = \tau_{\{\varepsilon\} \cup \text{succ}(P_d) \cup \text{succ}(v)} \}$ 

```

```

do  $u \neq \varepsilon \rightarrow$ 
  if  $\text{tau}(v, u \upharpoonright 1) = \perp \rightarrow \text{tau}(v, u \upharpoonright 1) := (u \upharpoonright 1)v;$ 
  for  $a : a \in V \rightarrow \text{tau}((u \upharpoonright 1)v, a) := \perp$  rof
   $\|\text{tau}(v, u \upharpoonright 1) \neq \perp \rightarrow \text{skip}$ 
  fi;
   $u, v := u \downarrow 1, (u \upharpoonright 1)v$ 
od;
 $P_d, P_r := P_d + \{p\}, P_r - \{p\}$ 
od $\{\text{tau} = \tau_P\}.$ 

```

In the algorithms we use $+$ for the union of disjoint sets and $-$ for the difference of a set and a subset of it. Notice that variable P_d is only needed to formulate an invariant for tau , so it may safely be removed from the algorithm. Furthermore, the states of the reverse trie are represented by strings. In practice, one can resort to a more suitable representation, for instance a representation by natural numbers. We will not elaborate this here.

4.2. Precomputation of d -functions

In this section we show that all d -functions introduced in Section 3 can be written according to a general pattern. For this general pattern a general precomputation algorithm is derived. In order to obtain a precomputation algorithm for a particular d -function one only has to instantiate the general precomputation algorithm and possibly simplify the resulting algorithm. The general pattern we strive for is a function $d \in V \times \text{su}ff(P) \rightarrow \mathbb{N}$ defined for $a \in V$ and $y \in \text{su}ff(P)$ by

$$d(a, y) = (\text{MIN } t : t \in \text{su}ff(P) \setminus \{\varepsilon\} \wedge Q(a, t) \wedge R(a, y, t) \wedge y <_p t : |t| - |y|)$$

where Q is a predicate on $V \times V^*$ and R a predicate on $V \times V^* \times V^*$. Why both Q and R are introduced will become clear when we derive an algorithm scheme for the computation of d . We will now show that all d -functions introduced in Section 3 can be expressed in this pattern. In the following derivations let $a \in V$ and $y \in \text{su}ff(P)$. First we derive

$$\begin{aligned}
& d_{opt}(a, y) \\
= & \{\text{definition } d_{opt}\} \\
& (\text{MIN } n : n \geq 1 \wedge V^* a y V^n \cap P \neq \emptyset : n) \\
= & \{\text{Property B.4}\} \\
& (\text{MIN } n : n \geq 1 \wedge a y V^n \cap \text{su}ff(P) \neq \emptyset : n) \\
= & \{\text{change of bound variable: } n = |s|\} \\
& (\text{MIN } s : s \in V^+ \wedge a y s \in \text{su}ff(P) : |s|) \\
= & \{\text{change of bound variable: } t = ys, t \neq \varepsilon\} \\
& (\text{MIN } t : t \in \text{su}ff(P) \setminus \{\varepsilon\} \wedge at \in \text{su}ff(P) \wedge y <_p t : |t| - |y|).
\end{aligned}$$

Hence, we have expressed function d_{opt} according to the general pattern with $Q(a, t) = at \in \text{su}ff(P)$ and $R(a, y, t) = \text{true}$. Notice that $at \in \text{su}ff(P) \equiv \tau(t, a) \neq \perp$. Assuming $y = \varepsilon$ we derive

$$\begin{aligned}
 & d_{sp}(\varepsilon) \\
 = & \{ \text{definition } d_{sp} \} \\
 & (\text{MIN } n : n \geq 1 \wedge V^n \cap V^*P \neq \emptyset : n) \\
 = & \{ \text{calculus} \} \\
 & (\text{MIN } t : t \in P : (\text{MIN } n : n \geq 1 \wedge V^n \cap V^*t \neq \emptyset : n)) \\
 = & \{ \text{range split: } P = P \setminus \{\varepsilon\} \cup (P \cap \{\varepsilon\}) \} \\
 & (\text{MIN } t : t \in P \setminus \{\varepsilon\} : |t|) \text{ min } (\text{MIN } t : t \in P \cap \{\varepsilon\} : 1) \\
 = & \{ \text{rewriting in order to obtain general pattern} \} \\
 & (\text{MIN } t : t \in \text{su}ff(P) \setminus \{\varepsilon\} \wedge t \in P \wedge \varepsilon <_p t : |t| - |\varepsilon|) \\
 & \text{min } (\text{MIN } t : t \in P \cap \{\varepsilon\} : 1).
 \end{aligned}$$

Assuming $y \neq \varepsilon$ we derive

$$\begin{aligned}
 & d_{sp}(y) \\
 = & \{ \text{definition } d_{sp} \} \\
 & (\text{MIN } n : n \geq 1 \wedge yV^n \cap V^*P \neq \emptyset : n) \\
 = & \{ \text{Property B.4} \} \\
 & (\text{MIN } n : n \geq 1 \wedge \text{su}ff(yV^n) \cap P \neq \emptyset : n) \\
 = & \{ y \neq \varepsilon, \text{ hence } \text{su}ff(yV^n) = yV^n \cup \text{su}ff((y \downarrow 1)V^n), \text{ set calculus} \} \\
 & (\text{MIN } n : n \geq 1 \wedge yV^n \cap P \neq \emptyset : n) \\
 & \text{min } (\text{MIN } n : n \geq 1 \wedge \text{su}ff((y \downarrow 1)V^n) \cap P \neq \emptyset : n) \\
 = & \{ \text{change of bound variable: } n = |s|, \text{ Property B.4} \} \\
 & (\text{MIN } s : s \in V^+ \wedge ys \in P : |s|) \\
 & \text{min } (\text{MIN } n : n \geq 1 \wedge (y \downarrow 1)V^n \cap V^*P \neq \emptyset : n) \\
 = & \left\{ \begin{array}{l} \text{change of bound variable: } t = ys, t \neq \varepsilon \\ y \in \text{su}ff(P), y \downarrow 1 \in \text{su}ff(P), \text{ definition } d_{sp} \end{array} \right\} \\
 & (\text{MIN } t : t \in P \setminus \{\varepsilon\} \wedge y <_p t : |t| - |y|) \text{ min } d_{sp}(y \downarrow 1) \\
 = & \{ \text{rewriting in order to obtain general pattern} \} \\
 & (\text{MIN } t : t \in \text{su}ff(P) \setminus \{\varepsilon\} \wedge t \in P \wedge y <_p t : |t| - |y|) \text{ min } d_{sp}(y \downarrow 1).
 \end{aligned}$$

Although the derived definition of d_{sp} is recursive and d_{sp} does not have an argument $a \in V$ one can still discern the general pattern with $Q(a, t) = t \in P$ and $R(a, y, t) = \text{true}$. Precomputation of d_{sp} can be done according to the general precomputation

algorithm without an iteration over V , followed by a breadth first traversal of the reverse trie. Next we derive

$$\begin{aligned}
 & d_i(y) \\
 = & \{ \text{definition } d_i \} \\
 & (\text{MIN } n : n \geq 1 \wedge V^* y V^n \cap P \neq \emptyset : n) \\
 = & \{ \text{Property B.4} \} \\
 & (\text{MIN } n : n \geq 1 \wedge y V^n \cap \text{succ}(P) \neq \emptyset : n) \\
 = & \{ \text{change of bound variable: } n = |s| \} \\
 & (\text{MIN } s : s \in V^+ \wedge y s \in \text{succ}(P) : |s|) \\
 = & \{ \text{change of bound variable: } t = ys, t \neq \varepsilon \} \\
 & (\text{MIN } t : t \in \text{succ}(P) \setminus \{ \varepsilon \} \wedge y <_p t : |t| - |y|).
 \end{aligned}$$

Although d_i does not have an argument $a \in V$ its definition still matches the general pattern with $Q(a, t) = \text{true}$ and $R(a, y, t) = \text{true}$. We conclude by deriving

$$\begin{aligned}
 & d_{vi}(y) \\
 = & \{ \text{definition } d_{vi} \} \\
 & (\text{MIN } n : n \geq 1 \wedge V^* (V \setminus MS(y)) y V^n \cap P \neq \emptyset : n) \\
 = & \{ \text{Property B.4} \} \\
 & (\text{MIN } n : n \geq 1 \wedge (V \setminus MS(y)) y V^n \cap \text{succ}(P) \neq \emptyset : n) \\
 = & \{ \text{change of bound variable: } n = |s| \} \\
 & (\text{MIN } s : s \in V^+ \wedge (V \setminus MS(y)) y s \cap \text{succ}(P) \neq \emptyset : |s|) \\
 = & \{ \text{change of bound variable: } t = ys, t \neq \varepsilon \} \\
 & (\text{MIN } t : t \in \text{succ}(P) \setminus \{ \varepsilon \} \\
 & \quad \wedge (V \setminus MS(y)) t \cap \text{succ}(P) \neq \emptyset \wedge y <_p t : |t| - |y|) \\
 = & \{ \text{definition } MS \} \\
 & (\text{MIN } t : t \in \text{succ}(P) \setminus \{ \varepsilon \} \\
 & \quad \wedge MS(t) \cap (V \setminus MS(y)) \neq \emptyset \wedge y <_p t : |t| - |y|).
 \end{aligned}$$

Apart from the fact that d_{vi} does not have an argument $a \in V$ its definition still resembles the general pattern if we take $Q(a, t) = \text{true}$ and $R(a, y, t) = MS(t) \cap (V \setminus MS(y)) \neq \emptyset$.

Having expressed all d -functions from Section 3 in the general pattern we proceed by giving a rather straightforward and non-deterministic algorithm to compute d which will serve as a starting point for further algorithm derivations (notice that program

variable *dee* is used to compute and store function *d*):

```

for  $y, a : y \in \text{su}ff(P) \wedge a \in V \rightarrow \text{dee}(a, y) := +\text{inf}$  rof;
for  $t, a : t \in \text{su}ff(P) \setminus \{\varepsilon\} \wedge a \in V \rightarrow$ 
  if  $Q(a, t) \rightarrow$ 
    for  $y : y \in \text{su}ff(P) \wedge y <_p t \rightarrow$ 
      if  $R(a, y, t) \rightarrow \text{dee}(a, y) := \text{dee}(a, y) \min(|t| - |y|)$ 
       $\square \neg R(a, y, t) \rightarrow \text{skip}$ 
    fi
  rof
 $\square \neg Q(a, t) \rightarrow \text{skip}$ 
fi
rof  $\{\text{dee} = d\}$ .

```

For d_i and d_{sp} the iterations over V can be omitted. For d_{sp} the following additional breadth first traversal of the reverse trie ($\text{su}ff(P)$) is needed to complete the computation:

```

 $n := 1;$ 
do  $\text{su}ff(P) \cap V^n \neq \emptyset \rightarrow$ 
  for  $t : t \in \text{su}ff(P) \cap V^n \rightarrow \text{deesp}(t) := \text{deesp}(t) \min \text{deesp}(t \downarrow 1)$  rof
od  $\{\text{deesp} = d_{sp}\}$ .

```

We now concentrate on making the innermost repetition in the general algorithm deterministic. We start by defining function $sp \in \text{su}ff(P) \rightarrow \mathcal{P}(\text{su}ff(P))$ for $t \in \text{su}ff(P)$ by

$$sp(t) = \{y \mid y \in \text{su}ff(P) \wedge y <_p t\},$$

the set of all suffixes of keywords that are a proper prefix to t . Notice that for all $t \in \text{su}ff(P) \setminus \{\varepsilon\}$ set $sp(t)$ is finite, non-empty, and linearly ordered with respect to \leq_p . Therefore, we can define $m_{sp} \in \text{su}ff(P) \setminus \{\varepsilon\} \rightarrow \text{su}ff(P)$ for $t \in \text{su}ff(P) \setminus \{\varepsilon\}$ by (the MAX_{\leq_p} quantification is described in Appendix A)

$$m_{sp}(t) = (\text{MAX}_{\leq_p} y : y \in \text{su}ff(P) \wedge y <_p t : y),$$

the maximal element of $sp(t)$. In [2, 3, 5, 14], function m_{sp} is known as the *failure function* corresponding to the reverse trie. For $t \in \text{su}ff(P) \setminus \{\varepsilon\}$ we derive a recursive definition of $sp(t)$ in terms of function m_{sp} :

$$\begin{aligned}
 & sp(t) \\
 = & \quad \{\text{definition } sp\}
 \end{aligned}$$

$$\begin{aligned}
& \{y \mid y \in \mathbf{suff}(P) \wedge y <_p t\} \\
= & \{t \in \mathbf{suff}(P) \setminus \{\varepsilon\}, \text{ Property B.6}\} \\
& \{y \mid y \in \mathbf{suff}(P) \wedge (y = msp(t) \vee y <_p msp(t))\} \\
= & \{\text{definition } sp \text{ and } msp\} \\
& \{msp(t)\} \cup sp(msp(t)).
\end{aligned}$$

Provided msp is already computed (precomputation of msp is discussed in Section 4.3) the innermost repetition traversing the set $sp(t)$ can be replaced by the following deterministic repetition (variable \bar{v} is a ghost variable needed to express the invariant):

```

v := t; {  $\bar{v} := \emptyset$ ; }
{invariant:  $v \in sp(t) \cup \{t\} \wedge sp(t) = \bar{v} \cup sp(v) \wedge \bar{v} \cap sp(v) = \emptyset$ }
do  $v \neq \varepsilon \rightarrow \{sp(v) = \{msp(v)\} \cup sp(msp(v))\}$ 
     $v := msp(v)$ ; {  $\bar{v} := \bar{v} + \{v\}$ ; }
    if  $R(a, v, t) \rightarrow dee(a, v) := dee(a, v) \min(|t| - |v|)$ 
     $\neg R(a, v, t) \rightarrow \text{skip}$ 
fi
od {  $v = \varepsilon$ , so  $sp(v) = \emptyset$  and  $sp(t) = \bar{v}$  }.

```

The invariant expresses that we have a bipartition of $sp(t)$ in \bar{v} (elements of $sp(t)$ that have already contributed to the computation of d) and $sp(v)$ (the other elements of $sp(t)$).

In case $R(a, y, t) = \text{true}$ for all a, y , and t the inner repetition can be made more efficient. Notice that this can be done for all presented d -functions except d_{vi} . In the following assume that $R(a, y, t) = \text{true}$ for all a, y , and t . Suppose that for some $v \in sp(t)$ in the above repetition we have $|t| - |v| \geq dee(a, v)$. From the structure of the algorithms we infer that $dee(a, v) = |t_0| - |v|$ for some $t_0 \in \mathbf{suff}(P) \setminus \{\varepsilon\}$ with $|t_0| \leq |t|$ that has already contributed to the computation of d . Therefore, for all $s \in sp(v)$ we have $dee(a, s) \leq |t_0| - |s|$ leading to

$$\begin{aligned}
& dee(a, s) \\
\leq & \{ \} \\
& |t_0| - |s| \\
= & \{dee(a, v) = |t_0| - |v|\} \\
& dee(a, v) + |v| - |s| \\
\leq & \{dee(a, v) \leq |t| - |v|\} \\
& |t| - |s|.
\end{aligned}$$

Hence, the contribution of t will not change the already computed value of $dee(a, s)$ for $s \in sp(v)$ and the inner repetition can be terminated. This yields the following repetition using an additional boolean variable *contributes* (notice that ghost variable \bar{v} is omitted):

```

v := t; contributes := true ;
do v ≠ ε ∧ contributes → v := msp(v);
  if |t| - |v| < dee(a, v) → dee(a, v) := |t| - |v|
  □ |t| - |v| ≥ dee(a, v) → contributes := false
fi
od.

```

Variable *contributes* can be removed resulting in the following repetition

```

v := t;
do v ≠ ε → v := msp(v);
  if |t| - |v| < dee(a, v) → dee(a, v) := |t| - |v|
  □ |t| - |v| ≥ dee(a, v) → v := ε
fi
od.

```

In order to further exploit this phenomenon the elements of $\text{su}ff(P) \setminus \{\varepsilon\}$ are dealt with in order of increasing length, i.e. the outermost repetition of the general precomputation algorithm does a breadth first traversal of the reverse trie. This results in the following algorithm:

```

for y, a : y ∈ suff(P) ∧ a ∈ V → dee(a, y) := +inf rof;
n := 1;
do suff(P) ∩ Vn ≠ ∅ →
  for t, a : t ∈ suff(P) ∩ Vn ∧ a ∈ V →
    if Q(a, t) → v := t;
      do v ≠ ε →
        v := msp(v);
        if |t| - |v| < dee(a, v) → dee(a, v) := |t| - |v|
        □ |t| - |v| ≥ dee(a, v) → v := ε
      fi
    od
  □ ¬Q(a, t) → skip
  fi
rof;
n := n + 1
od.

```

In this optimized breadth first precomputation algorithm for each node v in the reverse trie and each symbol from V the step from v to $msp(v)$ is done at most two times. Therefore, the precomputation time is $O(|\mathbf{suff}(P)| \cdot |V|)$. If the traversal of V can be omitted (as is the case for the precomputation of d_i and d_{sp}) it is $O(|\mathbf{suff}(P)|)$. The breadth first precomputation algorithm for d_i can be simplified further by observing that since $Q(a, t) = \text{true}$ the steps taken from t are always preceded by the steps taken from $msp(t)$ (provided $msp(t) \neq \varepsilon$) since $|msp(t)| < |t|$. So, only the contribution of t to $d_i(msp(t))$ has to be considered. This results in the following precomputation algorithm for d_i (program variable $deei$ is used to compute and finally store d_i):

```

for  $y : y \in \mathbf{suff}(P) \rightarrow deei(y) := +\text{inf}$  rof;
 $n := 1$ ;
do  $\mathbf{suff}(P) \cap V^n \neq \emptyset \rightarrow$ 
  for  $t : t \in \mathbf{suff}(P) \cap V^n \rightarrow$ 
     $deei(msp(t)) := deei(msp(t)) \mathbf{min}(|t| - |msp(t)|)$ 
  rof;
   $n := n + 1$ 
od.
```

Notice that the breadth first traversal in this algorithm may be replaced by an arbitrary traversal of the reverse trie.

4.3. Precomputation of msp

Finally, we derive an algorithm computing function msp defined in Section 4.2. For $a \in \mathbf{suff}(P) \cap V$ we have

$$\begin{aligned}
 & msp(a) \\
 = & \{ \text{definition } msp \} \\
 & (\mathbf{MAX}_{\leq_p} y : y \in \mathbf{suff}(P) \wedge y <_p a : y) \\
 = & \{ y <_p a \equiv y = \varepsilon, \varepsilon \in \mathbf{suff}(P) \} \\
 & \varepsilon
 \end{aligned}$$

and for $ax \in \mathbf{suff}(P)$, $a \in V$, and $x \in \mathbf{suff}(P) \setminus \{\varepsilon\}$

$$\begin{aligned}
 & msp(ax) \\
 = & \{ \text{definition } msp \} \\
 & (\mathbf{MAX}_{\leq_p} y : y \in \mathbf{suff}(P) \wedge y <_p ax : y) \\
 = & \{ \text{range split, } \varepsilon <_p ax, \varepsilon \in \mathbf{suff}(P) \} \\
 & (\mathbf{MAX}_{\leq_p} y : y \in \mathbf{suff}(P) \wedge y <_p ax \wedge y \neq \varepsilon : y) \mathbf{max}_{\leq_p} \varepsilon
 \end{aligned}$$

```

for  $a : a \in \text{su}ff(P) \cap V \rightarrow \text{em}sp(a) := \varepsilon$  rof;
 $n := 1$ ;
{invariant:  $(\forall y : y \in \text{su}ff(P) \setminus \{\varepsilon\} \wedge |y| \leq n : \text{em}sp(y) = \text{msp}(y))$ }
do  $\text{su}ff(P) \cap V^n \neq \emptyset \rightarrow$ 
  for  $t, a : t \in \text{su}ff(P) \cap V^n \wedge a \in V \rightarrow$ 
    if  $at \in \text{su}ff(P) \rightarrow$ 
       $v := \text{em}sp(t)$ ;
      {linear search}
      do  $av \notin \text{su}ff(P) \wedge v \neq \varepsilon \rightarrow v := \text{em}sp(v)$  od;
      if  $av \in \text{su}ff(P) \rightarrow \text{em}sp(at) := av$ 
       $\square av \notin \text{su}ff(P) \wedge v = \varepsilon \rightarrow \text{em}sp(at) := \varepsilon$ 
      fi
     $\square at \notin \text{su}ff(P) \rightarrow \text{skip}$ 
  fi
  rof;
   $n := n + 1$ 
od.

```

Fig. 4. Precomputation algorithm for msp .

$$\begin{aligned}
&= \{\text{change of bound variable: } y = ay'\} \\
&(\text{MAX}_{\leq_p} y' : ay' \in \text{su}ff(P) \wedge y' <_p x : ay') \text{max}_{\leq_p} \varepsilon \\
&= \left\{ \begin{array}{l} ay' \in \text{su}ff(P) \Rightarrow y' \in \text{su}ff(P), \\ x \in \text{su}ff(P) \setminus \{\varepsilon\}, \text{definition } sp \end{array} \right\} \\
&(\text{MAX}_{\leq_p} y' : ay' \in \text{su}ff(P) \wedge y' \in sp(x) : ay') \text{max}_{\leq_p} \varepsilon.
\end{aligned}$$

From this it follows that $\text{msp}(ax)$ can be computed by a linear search in downward order over $sp(x)$ (remember that $sp(x)$ is linearly ordered with respect to \leq_p) starting with $\text{msp}(x)$. Provided the computation of msp is done using a breadth first traversal of the reverse trie ($\text{su}ff(P)$) the value of msp is already computed for all elements of $sp(x) \cup \{x\}$ and can therefore be used to implement the linear search over $sp(x)$. This results in the algorithm in Fig. 4 where variable $\text{em}sp$ is used to compute and finally store msp . This breadth first algorithm computing msp can be combined with the breadth first algorithm that computes the d -functions. The precomputation time is $O(|\text{su}ff(P)| \cdot |V| \cdot (\text{MAX}_p : p \in P : |p|))$. The precomputation time can be reduced to $O(|\text{su}ff(P)| \cdot |V|)$ at the expense of $O(|\text{su}ff(P)| \cdot |V|)$ additional storage space by simultaneously computing and storing the transition function of the deterministic finite

automaton mentioned in [2]. Since introduction of that transition function is not straightforward in the setting of this article we refer to [18,19] for the details of this approach.

5. Conclusions

In this article we derived and presented a taxonomy of sublinear keyword pattern matching algorithms closely related to the Boyer–Moore algorithm [3] and the Commentz-Walter algorithm [5,6]. It includes, amongst others, the multiple keyword generalization of the single keyword Boyer–Moore algorithm and the algorithm presented by Fan and Su [9,10]. We presented the algorithms within a common framework permitting an easier comprehension of and a better comparison between the algorithms. This was achieved by the systematic and formal derivation of the algorithms from a common starting point and by factoring out common portions of the derivations. The derivations were done through series of refinements to either algorithm or problem. A refinement to the algorithm/problem is referred to as the introduction of an algorithm/problem detail. The sequence of details that are subsequently introduced in a derivation characterizes the algorithm obtained by that derivation. Detail sequences can therefore be used to classify the algorithms in the taxonomy. Algorithms can now be compared by looking at their detail sequences. The taxonomy graph in Fig.1 constitutes a concise presentation and classification of the pattern matching algorithms discussed, vertices representing algorithms and edges representing the addition of an algorithm or problem detail. It can be viewed as a table of contents to this article. Our results show how fruitful the applied method of developing a taxonomy is (it was inspired by the method described by Jonkers [13]).

Introduction of the notion of safe shift distances proved to be essential for the derivation of the various algorithms. All algorithms are characterized by a – systematically derived and more or less easy to compute – approximation from below of the maximal safe shift distance, computation of the latter being equivalent to the keyword pattern matching problem itself. The systematic derivation provided a means to compare the algorithms and their matching speeds, and to get a better understanding of the algorithms and their interrelations. Perhaps this better understanding will help further the use of the algorithms from this family. Our derivations show the Commentz-Walter algorithm not to be the multiple keyword generalization of the Boyer–Moore algorithm (as was the original intention of Commentz-Walter) and that such a generalization can indeed be obtained. Of the algorithms presented the algorithm by Fan and Su [9,10] is the fastest (at the expense of additional precomputation time and additional storage requirements), followed by the common ancestor of the Boyer–Moore and Commentz-Walter algorithms, and then by both the multiple keyword generalization of a Boyer–Moore algorithm [3] and the Commentz-Walter algorithm [5,6]. The latter two are incomparable in matching speed. It is clear that we have not derived and presented all possible sublinear pattern matching algorithms. Our derivation method, however, clearly indicates how yet other members of this family of algorithms may be derived.

Apart from giving a taxonomy of sublinear pattern matching algorithms we presented the first formally derived and therefore correct precomputation algorithms (this cannot always be said about the algorithms found in the literature, mostly due to the absence of any formal derivation; see, for instance, the many solutions for the Boyer–Moore precomputation that have been published, corrected and republished). In fact, we showed that most of the precomputation algorithms can be obtained as instantiations of a general precomputation algorithm scheme derived for a general function pattern in which most components of the various shift functions can be expressed. Thus, we provided a common framework for the precomputation algorithms as well.

Acknowledgements

We thank Kees Hemerik, Frans Kruseman Aretz, Anne Kaldewaij, and Lex Bijlsma for their careful reading of earlier drafts of this article and their constructive criticisms. We also thank the two referees for their valuable comments.

Appendix A. Quantifications

In this article we use a non-standard notation for quantifications more suited for the calculational style in which we derive algorithms. Let X be a set and \oplus a commutative, associative, and idempotent operator on X (i.e. $\oplus : X \times X \rightarrow X$) having unity e . Let Y be some other set, and $E(y)$ an expression of type X for each $y \in Y$. Let P be a predicate on Y such that $\{y \in Y \mid P(y)\}$ is finite, say $\{y \in Y \mid P(y)\} = \{y_1, y_2, \dots, y_n\}$ for some n . The expression

$$E(y_1) \oplus E(y_2) \oplus \dots \oplus E(y_n)$$

will be written as

$$(\oplus y : P(y) : E(y)).$$

This is called a *quantification*, \oplus the *quantifier*, y the *bound variable*, P the *range* (*predicate*), and $E(y)$ the *term*. Furthermore, we define

$$(\oplus y : \text{false} : E(y)) = e \quad (\text{the unity of } \oplus).$$

For instance, taking $X = \mathcal{P}(\mathbb{N})$ (the powerset of \mathbb{N}), $\oplus = \bigcup$, $e = \emptyset$, $Y = \mathbb{N}$, $E(y) = \{y\}$, and $P(y) \equiv 0 \leq y < N$ then

$$(\bigcup y : 0 \leq y < N : \{y\})$$

denotes the set $\{0, 1, \dots, N-1\}$ or

$$\bigcup_{y: 0 \leq y < N} \{y\}$$

Table 1

\oplus	\bigoplus	e
\wedge	\forall	true
\vee	\exists	false
\cup	\bigcup	\emptyset
min	MIN	$+\infty$
\max_{\leq_p}	MAX_{\leq_p}	ε

in a more traditional notation, and

$$(\bigoplus y : 0 \leq y < 0 : \{y\}) = \emptyset.$$

Table 1 lists a number of operators, their corresponding quantifiers, and their unities.

The last operator is defined on a set of strings that is linearly ordered with respect to \leq_p (see Definition B.3).

Property A.1. *Quantifications satisfy*

$$(\bigoplus y : P_1(y) \vee P_2(y) : E(y)) = (\bigoplus y : P_1(y) : E(y)) \oplus (\bigoplus y : P_2(y) : E(y))$$

(described by hint range split or disjunctive range).

Property A.2. *The MIN-quantifier satisfies*

- (i) if $P \Rightarrow Q$ then $(\text{MIN } y : P(y) : E(y)) \geq (\text{MIN } y : Q(y) : y)$ (described by hint enlarging range),
- (ii) $(\text{MIN } y : P(y) \wedge Q(y) : E(y)) \geq (\text{MIN } y : P(y) : E(y)) \text{ max } (\text{MIN } y : Q(y) : E(y))$ (described by hint conjunctive range),
- (iii) $(\text{MIN } y : P(y) : E(y) + c) = (\text{MIN } y : P(y) : E(y)) + c$ provided $(\exists y : \text{true} : P(y))$ (described by hint non-empty range).

Appendix B. Definitions and properties

This section provides a series of definitions and properties which are used throughout this article. In the following let V be an alphabet.

For a string $w \in V^*$ w^R denotes the reversal of w . For any language $L \subseteq V^*$ we define $L^R = \{w^R \mid w \in L\}$ (the reversal of language L).

Definition B.1. The infix operators $\uparrow, \downarrow, \upharpoonright, \downharpoonright : V^* \times \mathbb{N} \rightarrow V^*$ are defined for $a \in V$, $v, w \in V^*$ and $k \in \mathbb{N}$ by

$$v \uparrow 0 = \varepsilon,$$

$$\varepsilon \uparrow (k + 1) = \varepsilon,$$

$$\begin{aligned}
(aw) \upharpoonright (k+1) &= a(w \upharpoonright k), \\
v \downharpoonright 0 &= v, \\
\varepsilon \downharpoonright (k+1) &= \varepsilon, \\
(aw) \downharpoonright (k+1) &= w \downharpoonright k, \\
v \upharpoonright k &= (v^R \upharpoonright k)^R, \\
v \downharpoonright k &= (v^R \downharpoonright k)^R.
\end{aligned}$$

The operators \upharpoonright , \downharpoonright , \upharpoonright , and \downharpoonright are called *left take*, *left drop*, *right take*, and *right drop*, respectively.

Definition B.2. Functions $\mathbf{pref} : \mathcal{P}(V^*) \rightarrow \mathcal{P}(V^*)$ and $\mathbf{suff} : \mathcal{P}(V^*) \rightarrow \mathcal{P}(V^*)$ are defined for $L \subseteq V^*$ by

$$\mathbf{pref}(L) = \{w \mid w \in V^* \wedge (\exists x : x \in V^* : wx \in L)\}$$

and

$$\mathbf{suff}(L) = (\mathbf{pref}(L^R))^R = \{w \mid w \in V^* \wedge (\exists x : x \in V^* : xw \in L)\}.$$

For $w \in V^*$ we will write $\mathbf{pref}(w)$ ($\mathbf{suff}(w)$) instead of $\mathbf{pref}(\{w\})$ ($\mathbf{suff}(\{w\})$).

Definition B.3. The binary relations \leq_p and \leq_s over $V^* \times V^*$ are defined by $u \leq_p v \equiv u \in \mathbf{pref}(v)$ and $u \leq_s v \equiv u \in \mathbf{suff}(v)$.

The following two properties are used in the derivation of the Commentz-Walter precomputation algorithm.

Property B.4. Let $A, B \subseteq V^*$. Then $\mathbf{pref}(A) \cap B \neq \emptyset \equiv A \cap BV^* \neq \emptyset$ and $\mathbf{suff}(A) \cap B \neq \emptyset \equiv A \cap V^*B \neq \emptyset$.

Property B.5. Let $A, B \subseteq V^*$ and $a \in V$. Then

- (i) $V^*A \cap V^*B \neq \emptyset \equiv V^*A \cap B \neq \emptyset \vee A \cap V^*B \neq \emptyset$,
- (ii) $V^*aA \cap V^*B \neq \emptyset \equiv V^*aA \cap B \neq \emptyset \vee A \cap V^*B \neq \emptyset$,
- (iii) $V^*A \cap V^*B \neq \emptyset \equiv V^*A \cap B \neq \emptyset \vee A \cap V^+B \neq \emptyset$.

We conclude with a property of function \mathbf{msp} defined in Section 4.2.

Property B.6. For $x, y \in \mathbf{suff}(P)$ and $y \neq \varepsilon$ we have

$$x <_p y \equiv x \leq_p \mathbf{msp}(y).$$

Proof. Let $x, y \in \text{suff}(P)$ and $y \neq \varepsilon$. We derive

$$\begin{aligned}
 & x <_p y \\
 = & \{ \text{definition of } <_p \text{ and } \mathbf{pref}, x \in \text{suff}(P) \} \\
 & x \in (\mathbf{pref}(y) \setminus \{y\}) \cap \text{suff}(P) \\
 \Rightarrow & \left\{ \begin{array}{l} (\mathbf{pref}(y) \setminus \{y\}) \cap \text{suff}(P) \text{ is non-empty } (y \neq \varepsilon), \\ \text{finite and linearly ordered w.r.t. } \leq_p \end{array} \right\} \\
 & x \leq_p (\mathbf{MAX}_{\leq_p} w : w \in (\mathbf{pref}(y) \setminus \{y\}) \cap \text{suff}(P) : w) \\
 = & \{y \neq \varepsilon, \text{definition of } msp\} \\
 & x \leq_p msp(y) \\
 \Rightarrow & \{y \neq \varepsilon, msp(y) <_p y \text{ (by definition of } msp), \text{transitivity of } <_p\} \\
 & x <_p y. \quad \square
 \end{aligned}$$

Appendix C. Algorithm and problem details

In this appendix we list the algorithm and problem details introduced in this article with a short description.

P	examine prefixes of a given string in any order
S	examine suffixes of a given string in any order
+	examine the strings from a given set in order of increasing length (this program detail can only be applied after, for instance, program details P and S)
RT	usage of the transition function of the reverse trie corresponding to the set of keywords to check whether a string, that is a suffix of some keyword, preceded by a symbol is again a suffix of some keyword
SSD	allow any shift distance at least one that is safe, i.e. that does not cause the omission of any matches
NLAU	no lookahead at the symbols of the unscanned part of the input string when computing a safe shift distance
OLAU	one symbol lookahead at the unscanned part of the input string when computing a safe shift distance
OPT	When computing a safe shift distance use the recognized suffix and only the immediately preceding (mismatching) symbol, strictly coupled
NLA	when computing a safe shift distance do not look at the symbols preceding the recognized suffix
BMCW	when computing a safe shift distance, on the one hand, use the recognized suffix and the fact that the symbol preceding it is mismatching, and on the other hand, but strictly independent, the identity of that symbol
BM	lessen the contribution of the symbol preceding the recognized suffix to the shift distance in case it does not occur in any keyword

- cw when computing a safe shift distance do not use the fact that the symbol preceding the recognized suffix is mismatching (use the recognized suffix and the symbol preceding it independently)
- OKW the set of keywords contains only one keyword (in contrast to the preceding program details this is a problem detail).

References

- [1] A.V. Aho, Algorithms of finding patterns in strings, in: J. van Leeuwen, ed., *Handbook of Theoretical Computer Science*, Vol. A (Elsevier, Amsterdam, 1990) 255–300.
- [2] A.V. Aho and M.J. Corasick, Efficient string matching: An aid to bibliographic search, *Comm. ACM* **18** (1975) 333–340.
- [3] R.S. Boyer and J.S. Moore, A fast string searching algorithm, *Comm. ACM* **20** (1977) 762–772.
- [4] M. Broy, Program construction by transformations: a family tree of sorting programs, in: A.W. Biermann and G. Guihu, eds., *Computer Program Synthesis Methodologies*, NATO ASI, 1981, NATO ASI Series, Series C: Mathematical and Physical Sciences, Vol. 95 (Reidel, Dordrecht, 1983) 1–49.
- [5] B. Commentz-Walter, A string matching algorithm fast on the average, in: *Proc. ICALP '79*, Lecture Notes in Computer Science, Vol. 71 (Springer, Berlin, 1979) 118–132.
- [6] B. Commentz-Walter, A string matching algorithm fast on the average. Tech. Report TR 79.09.007, IBM-Germany, Scientific Center Heidelberg, 1979.
- [7] J. Darlington, A synthesis of several sorting algorithms, *Acta Inform.* **11** (1978) 1–30.
- [8] E.W. Dijkstra, *A Discipline of Programming* (Prentice-Hall, Englewood Cliffs, NJ, 1976).
- [9] J.-J. Fan and K.-Y. Su, An efficient algorithm for matching multiple patterns, *IEEE Trans. Knowledge Data Eng.* **5** (1993) 339–351.
- [10] J.-J. Fan and K.-Y. Su, An efficient algorithm for matching multiple patterns, in: J. Ae, ed., *Computer Algorithms: String Pattern Matching Strategies* (IEEE Computer Society Press, Los Alamitos, CA, 1994).
- [11] E. Fredkin, Trie memory, *Comm. ACM* **3** (1960) 490–499.
- [12] A. Hume and D.M. Sunday, Fast string searching, *Softw. Pract. Exper.* **21** (1991) 1221–1248.
- [13] H.B.M. Jonkers, *Abstraction, Specification and Implementation Techniques, with an Application to Garbage Collection*, Mathematical Centre Tracts, Vol. 166 (Mathematisch Centrum, Amsterdam, 1983).
- [14] D.E. Knuth, J.H. Morris, Jr and V.R. Pratt, Fast pattern matching in strings, *SIAM J. Comput.* **6** (1977) 323–350.
- [15] A.J.J.M. Marcelis, On the classification of attribute evaluation algorithms, *Sci. Comput. Programming* **14** (1) (1990) 1–24.
- [16] W. Rytter, A correct preprocessing algorithm for Boyer–Moore string-searching, *SIAM J. Comput.* **9** (1980) 509–512.
- [17] B.W. Watson, The performance of single-keyword and multiple-keyword pattern matching algorithms, Computing Science Report 94/19, Department of Computing Science, Eindhoven University of Technology, 1994.
- [18] B.W. Watson, Taxonomies and toolkits of regular language algorithms, Ph.D. Thesis, Department of Computing Science, Eindhoven University of Technology, 1995.
- [19] B.W. Watson and G. Zwaan, A taxonomy of keyword pattern matching algorithms. Computing Science Report 92/27, Department of Computing Science, Eindhoven University of Technology, 1992.
- [20] B.W. Watson and G. Zwaan, A taxonomy of keyword pattern matching algorithms, in: H. Wijshoff, ed., *Proc. Computing Science in the Netherlands '93* (Stichting Mathematisch Centrum, Amsterdam, 1993) 25–39.